

Orthogonal Sparse Eigenvectors: A Procrustes Problem

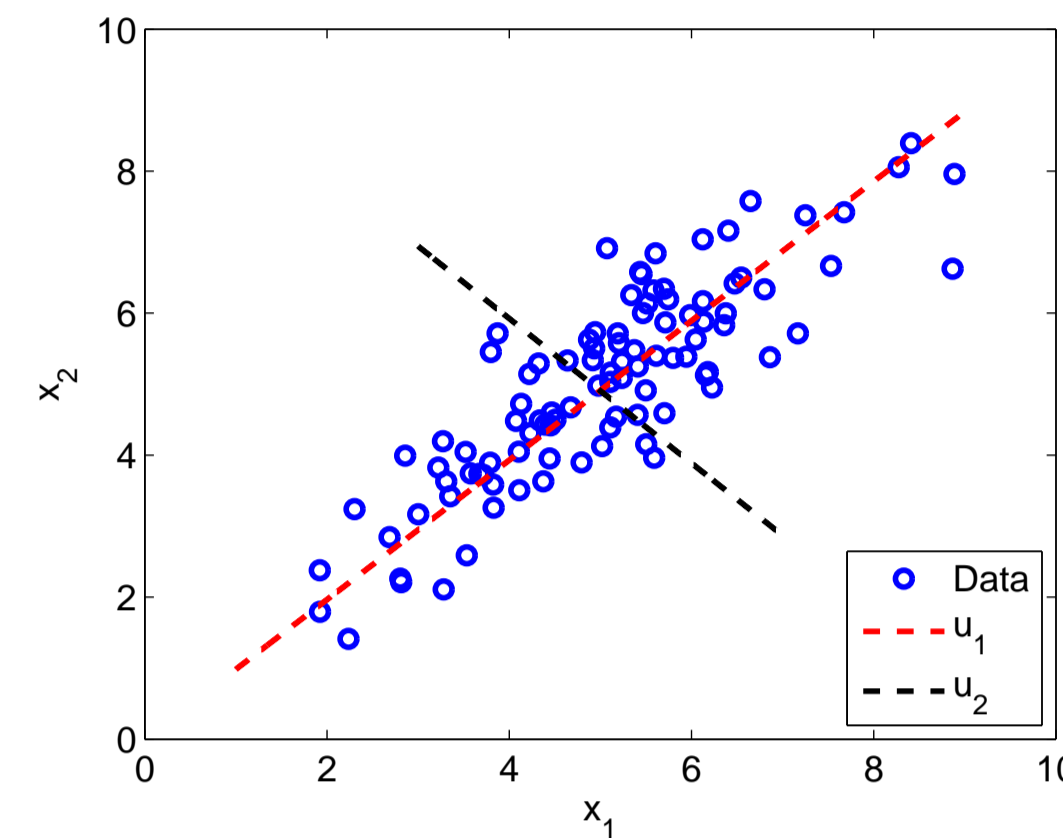
Konstantinos Benidis, Ying Sun, Prabhu Babu, Daniel P. Palomar

Department of Electronic and Computer Engineering, The Hong Kong University of Science and Technology, Hong Kong

Background & Motivation

- Principal Component Analysis (PCA) is a popular technique for data analysis and dimensionality reduction.
 - Captures directions of maximum variance of the data.
 - These directions (eigenvectors - PC loadings) form an orthonormal basis.
 - Principal components (PCs) are uncorrelated.

- Principal components are, in general, combinations of all the input variables.
 - PC loadings are dense vectors.
 - In many applications the variables have a physical meaning (e.g. gene expression).
 - A sparse basis would help significantly the interpretability of the result.
- Trade-offs:
 - Explained variance.
 - Orthogonality of the PC loadings.
 - Uncorrelatedness of the PCs.



Related Work

- Existing methods:
 - All the existing algorithms sacrifice orthogonality for a sparse result.
 - Benchmark: GPower (Journée et al. [2010]).
- Goal: Extract sparse eigenvectors that preserve the orthogonality property.**

Problem Formulation

- The orthogonal sparse eigenvector extraction translates to the following optimization problem:

$$\begin{aligned} & \underset{U}{\text{maximize}} \quad \text{Tr}(U^T S U D) - \sum_{i=1}^q \rho_i \|u_i\|_0 \\ & \text{subject to} \quad U^T U = I_q, \end{aligned}$$

where $U \in \mathbf{R}^{m \times q}$ denotes the eigenvectors, $S \in \mathbf{R}^{m \times m}$ the sample covariance matrix and $\|u_i\|_0$ the number of nonzero elements of the i -th eigenvector. $D = \text{Diag}(\mathbf{d}) \in \mathbf{R}_+^{q \times q}$ and ρ_i are regularization parameters.

- Without the sparsity (red) term it is the typical eigenvector extraction problem.
- Discontinuous, non-differentiable, non-concave objective function.
- Non-convex set.

Approximate Smooth Formulation

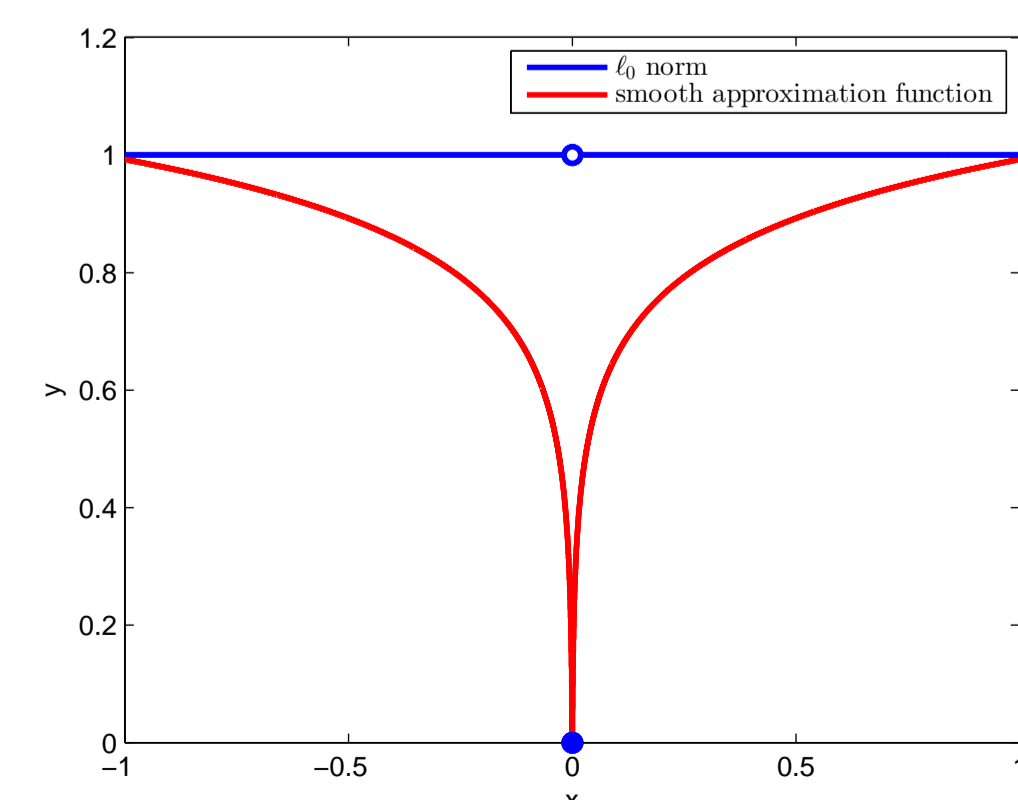
- We approximate the ℓ_0 norm with a smooth continuous and differentiable function (Song et al. [2015]):

$$\begin{aligned} & \underset{U}{\text{maximize}} \quad \text{Tr}(U^T S U D) - \sum_{j=1}^q \rho_j \sum_{i=1}^m g_p^\epsilon(u_{ij}) \\ & \text{subject to} \quad U^T U = I_q, \end{aligned} \quad (1)$$

where

$$g_p^\epsilon(x) = \begin{cases} \frac{x^2}{2\epsilon(p+\epsilon)\log(1+1/p)}, & |x| \leq \epsilon, \\ \frac{\log\left(\frac{p+|x|}{p+\epsilon}\right) + \frac{\epsilon}{2(p+\epsilon)}}{\log(1+1/p)}, & |x| > \epsilon, \end{cases}$$

- with $0 < p \leq 1$ and $0 < \epsilon \ll 1$.
- The problem is still non-convex.
 - Use the MM framework.



Minorization-Maximization Framework

- Problem:

$$\begin{aligned} & \underset{\mathbf{x}}{\text{maximize}} \quad f(\mathbf{x}) \\ & \text{subject to} \quad \mathbf{x} \in \mathcal{X}. \end{aligned}$$

- Minorization-Maximization algorithm:

$$\mathbf{x}^{(k+1)} = \arg \max_{\mathbf{x} \in \mathcal{X}} g(\mathbf{x} | \mathbf{x}^{(k)}).$$

- Surrogate function $g(\mathbf{x} | \mathbf{x}^{(k)})$ satisfies:

$$\begin{aligned} & f(\mathbf{x}^{(k)}) = g(\mathbf{x}^{(k)} | \mathbf{x}^{(k)}), \\ & f(\mathbf{x}) \geq g(\mathbf{x} | \mathbf{x}^{(k)}) \quad \forall \mathbf{x} \in \mathcal{X}, \\ & f'(\mathbf{x}^{(k)}; \mathbf{d}) = g'(\mathbf{x}^{(k)}; \mathbf{d} | \mathbf{x}^{(k)}), \\ & \forall \mathbf{x}^{(k)} + \mathbf{d} \in \mathcal{X}. \end{aligned}$$

- Iteratively maximize $g(\mathbf{x} | \mathbf{x}^{(k)})$ instead of maximizing $f(\mathbf{x})$.

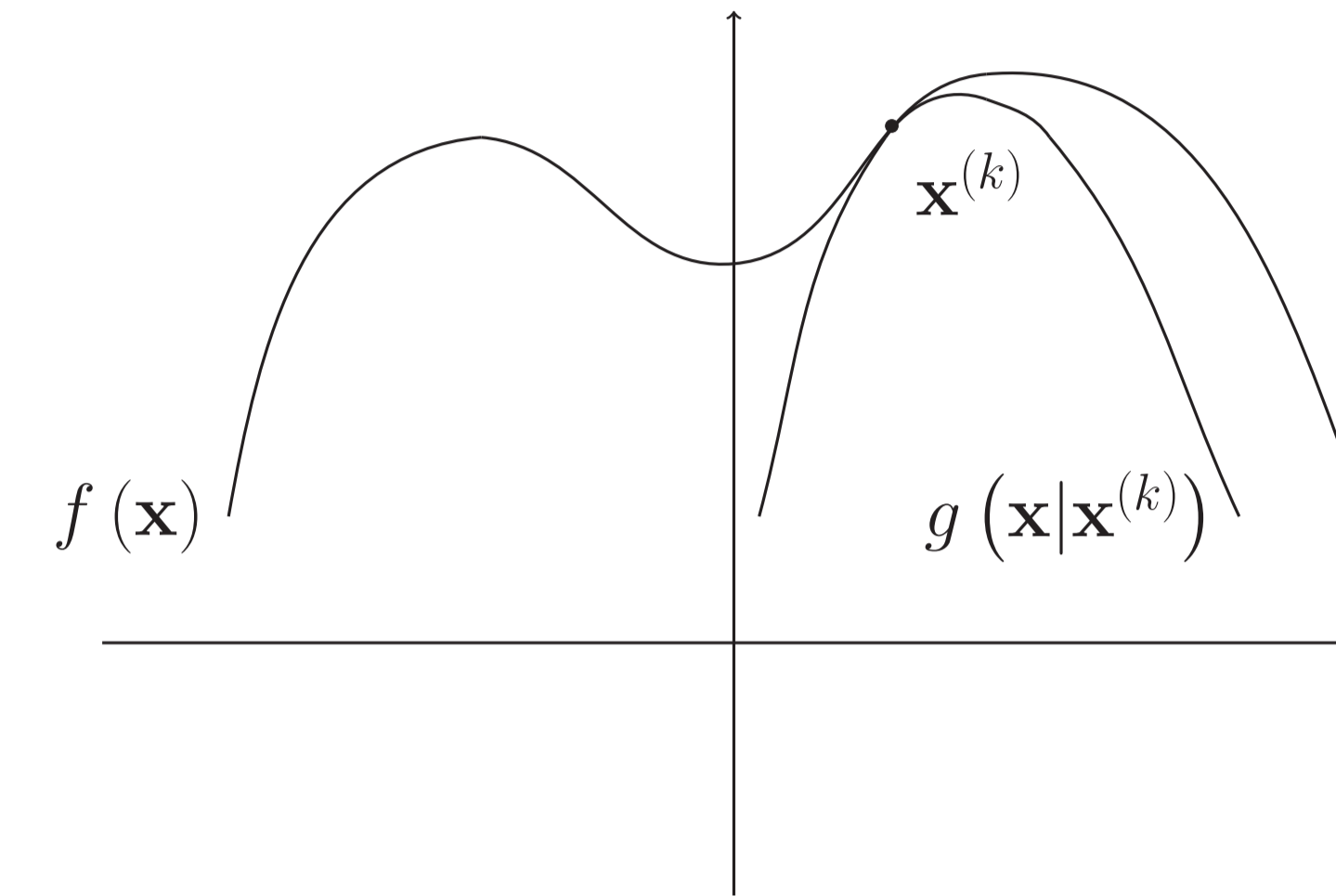


Figure: Minorization-Maximization Algorithm

Proposition

The objective function of (1) is lowerbounded by the surrogate function

$$g(U | U^{(k)}) = 2\text{Tr}\left(\left(G^{(k)} - H^{(k)}\right)^T U\right) + c,$$

where

$$G^{(k)} = S U^{(k)} D, \quad (2)$$

$$H^{(k)} = \left[\text{diag}\left(\mathbf{w}^{(k)} - \mathbf{w}_{\max}^{(k)} \otimes \mathbf{1}_m\right) \text{vec}\left(U^{(k)}\right) \right]_{m \times q}, \quad (3)$$

and c is an optimization irrelevant constants. The weights $\mathbf{w}^{(k)} \in \mathbf{R}_+^{mq}$ are given by

$$w_i^{(k)} = \begin{cases} \frac{\rho_i}{2\epsilon(p+\epsilon)\log(1+1/p)}, & |u_i^{(k)}| \leq \epsilon, \\ \frac{\rho_i}{2\log(1+1/p)|u_i^{(k)}|(|u_i^{(k)}|+p)}, & |u_i^{(k)}| > \epsilon, \end{cases}$$

where $\mathbf{u}^{(k)} = \text{vec}(U^{(k)})$, and $\mathbf{w}_{\max}^{(k)} \in \mathbf{R}_+^q$, with $w_{\max,i}^{(k)}$ being the maximum weight that corresponds to the i -th eigenvector. Equality is achieved when $U = U^{(k)}$.

Procrustes Reformulation

- Observe that

$$\arg \max_{U \in V_{m,q}} \text{Tr}\left(\left(G^{(k)} - H^{(k)}\right)^T U\right) = \arg \min_{U \in V_{m,q}} \|U - \left(G^{(k)} - H^{(k)}\right)\|_F^2$$

where $V_{m,q} = \{U \in \mathbf{R}^{m \times q} | U^T U = I_q\}$ denotes the Stiefel manifold.

- The optimization problem of every MM iteration takes the following form:

$$\begin{aligned} & \underset{U}{\text{minimize}} \quad \|U - \left(G^{(k)} - H^{(k)}\right)\|_F^2 \\ & \text{subject to} \quad U^T U = I_q. \end{aligned} \quad (4)$$

- The optimization problem (4) is a rectangular Procrustes problem.
 - Closed-form solution.

Lemma: Rectangular Procrustes

An optimal solution of the optimization problem (4) is $U^* = V_L V_R^T$, where V_L, V_R are the left and right singular vectors of the matrix $\left(G^{(k)} - H^{(k)}\right)$, respectively (Manton [2002]).

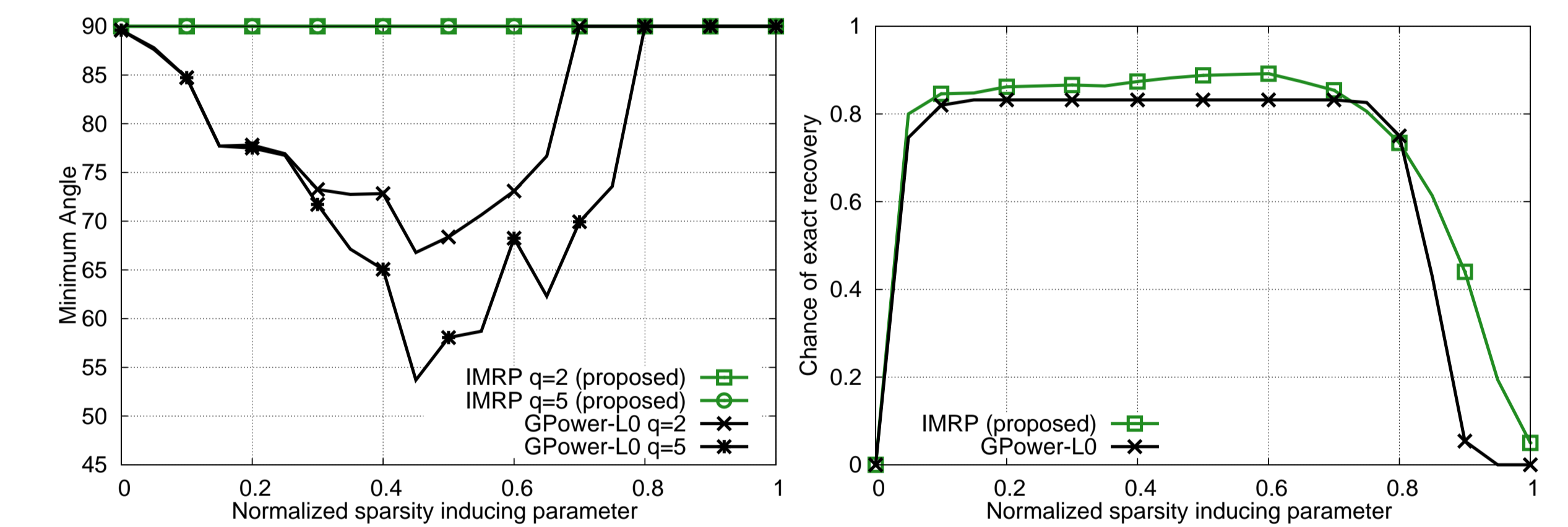
Algorithm

Algorithm 1 IMRP - Iterative Minimization of Rectangular Procrustes

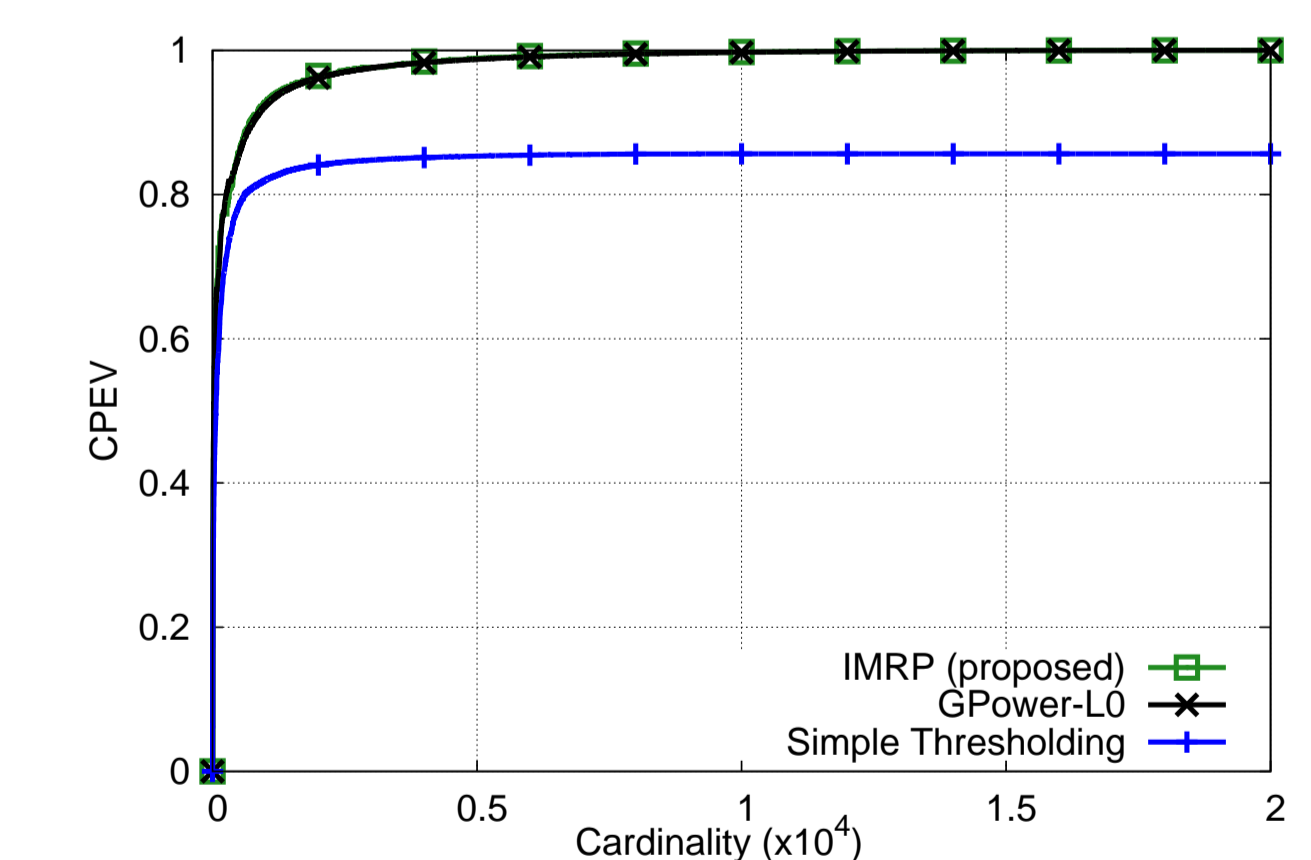
- Set $k = 0$, choose $U^{(0)} \in \{U : U^T U = I_q\}$
- repeat:**
- Compute $G^{(k)}, H^{(k)}$ with (2),(3), respectively
- Compute V_L, V_R , the left and right singular vectors of $\left(G^{(k)} - H^{(k)}\right)$, respectively
- $U^{(k+1)} = V_L V_R^T$
- $k \leftarrow k + 1$
- until** convergence
- return** $U^{(k)}$

Numerical Results

- Construct a covariance matrix Σ through the eigenvalue decomposition $\Sigma = V \text{diag}(\lambda) V^T$.
- The first q eigenvectors have a pre-specified sparse structure.
- We consider a setup with $m = 500, q = 5$.
- Generate 500 data matrices $A \in \mathbf{R}^{m \times n}$ by drawing $n = 50$ samples from a zero-mean normal distribution with covariance matrix Σ , i.e., $\mathbf{a}_i \sim \mathcal{N}(\mathbf{0}, \Sigma)$, for $i = 1, \dots, n$.



- Gene expression dataset: We consider $m = 4000$ genes with the largest variance and we estimate $q = 5$ sparse eigenvectors.



Conclusion

- We have proposed a new algorithm (IMRP) for sparse eigenvalue extraction.
 - Unlike all the other existing methods, the resulting sparse eigenvectors preserve the orthogonality property.
 - IMRP improves the chance of exact recovery of the eigenvectors and matches the cumulative percentage of explained variance (CPEV).

References

- Michel Journée, Yurii Nesterov, Peter Richtárik, and Rodolphe Sepulchre. Generalized power method for sparse principal component analysis. *The Journal of Machine Learning Research*, 11:517–553, March 2010.
- Junxiao Song, Prabhu Babu, and Daniel P. Palomar. Sparse generalized eigenvalue problem via smooth optimization. *IEEE Transactions on Signal Processing*, 63(7):1627–1642, April 2015.
- Jonathan H Manton. Optimization algorithms exploiting unitary constraints. *IEEE Transactions on Signal Processing*, 50(3):635–650, March 2002.